

Bi-convex Implicit Deep Learning

Bertrand Travacca

October 2019

1 Introduction

1.1 Implicit deep learning

Given an input $u \in \mathbb{R}^n$, where n denotes the number of features, as in [1], we define the *implicit deep learning* prediction rule $\hat{y}(u) \in \mathbb{R}^p$ with ReLU activation

$$\begin{aligned}\hat{y}(u) &= Ax + Bu + c \\ x &= (Dx + Eu + f)_+\end{aligned}\tag{1}$$

where $(\cdot)_+ := \max(0, \cdot)$ is the ReLU activation, $x \in \mathbb{R}^h$ is called the *hidden variable* (h is the number of hidden features), $\Theta := (A, B, c, D, E, f)$ are matrices and vectors of appropriate size, they define the parameters of the model. The hidden variable x is implicit in the sense that there is in general no analytical formula for it, this is different from classic deep learning for which, given the model parameters, the hidden variables can be computed explicitly via propagation through the network.

1.2 Notation and definitions

We denote $\|\cdot\|$ the euclidean norm, $\|\cdot\|_2$ the corresponding operator norm (i.e. the spectral norm) and $\|\cdot\|_F$ the Frobenius norm. \mathbb{R}_+^n denotes the positive orthant of the vector space \mathbb{R}^n , \mathbb{S}^n the set of real symmetric matrices of size n and \mathbb{S}_+^n the cone of positive semi-definite matrices of size n . The transpose of a matrix or vector is denoted $^\top$ and elementwise product is denoted \odot . Given a differentiable function f from $\mathbb{R}^{n \times p}$ to \mathbb{R} we define the scalar by matrix partial derivative in denominator layout convention as

$$\frac{\partial f}{\partial A} = \nabla_A f = \begin{bmatrix} \frac{\partial f}{\partial A_{1,1}} & \cdots & \frac{\partial f}{\partial A_{1,p}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{n,1}} & \cdots & \frac{\partial f}{\partial A_{n,p}} \end{bmatrix} \in \mathbb{R}^{n \times p}$$

We say that a function $(x, y) \rightarrow f(x, y)$ with separable domain of definition $\mathcal{X} \times \mathcal{Y}$ is bi-convex in (x, y) , if for all $x \in \mathcal{X}$, the function $y \rightarrow f(x, y)$ is convex and for all $y \in \mathcal{Y}$ the function $x \rightarrow f(x, y)$ is convex. We say that a function is smooth if it is differentiable and its gradient is Lipschitz continuous. We say that f is bi-smooth if it is smooth in x given y and vice-versa. An example of bi-smooth and bi-convex function is $(x, A) \rightarrow x^\top Ax$, $A \in \mathbb{S}_+^n$.

2 Well-posedness

We say that matrix D is well-posed for (1) if there exists a unique solution $x = (Dx + \delta)_+$, $\forall \delta \in \mathbb{R}^h$. Using the fact that ReLU is 1-Lipschitz we have for $x_1, x_2 \in \mathbb{R}^h$

$$\|(Dx_1 + \delta)_+ - (Dx_2 + \delta)_+\| \leq \|D(x_1 - x_2)\| \leq \|D\|_2 \|x_1 - x_2\|$$

If $\|D\|_2 < 1$ we have that the map $x \rightarrow (Dx + \delta)_+$ is a strict contraction. In that case, Banach's contraction mapping theorem applies, showing that the equation $x = (Dx + \delta)_+$ has a unique solution. In that case, a solution x can be computed via the Picard iterations

$$x^{k+1} = (Dx^k + \delta), \quad k = 1, 2, \dots$$

Note that $\|D\|_2 < 1$ is only a sufficient condition for well-posedness. Nevertheless this is the only condition we will consider in this article.

3 Learning with the implicit prediction rule

3.1 Problem formulation

Let us consider the input and output data matrices $U = [u_1, \dots, u_m] \in \mathbb{R}^{n \times m}$, $Y = [y_1, \dots, y_m] \in \mathbb{R}^{p \times m}$ – with m being the number of datapoints – the corresponding optimization regression problem (i.e. with squared error loss) reads

$$\begin{aligned} \min_{X, \Theta} \quad & \mathcal{L}(Y, [\Theta, X]) := \frac{1}{2m} \|AX + BU + c1_m^\top - Y\|_F^2 \\ \text{s.to:} \quad & X = (DX + EU + f1_m^\top)_+, \quad \|D\|_2 < 1 \end{aligned} \quad (2)$$

where 1_m is a column vector of size m consisting of ones. For clarity we have highlighted in blue the optimization variables. The non-convexity of this problem arises from the nonlinear implicit constraint and the matrix product terms AX and DX . In practice we replace the constraint $\|D\|_2 < 1$ by the closed convex form $\|D\|_2 \leq 1 - \varepsilon$, where $\varepsilon > 0$ small.

Remark 1. This constraint can be formulated as a linear matrix inequality (LMI) via Schur complement

$$\|D\|_2 \leq 1 - \varepsilon \iff \begin{bmatrix} (1 - \varepsilon)I_h & D \\ D^\top & (1 - \varepsilon)I_h \end{bmatrix} \in \mathbb{S}_+^{2h}$$

3.2 Fenchel Divergence Lagrangian Relaxation

Using Fenchel-Young inequality, it can be shown that the equation $x = (Dx + Eu + f)_+$ is equivalent to

$$\begin{cases} \mathcal{F}(x, Ax + Bu + c) = 0 \\ x \geq 0 \end{cases} \quad (3)$$

with the *Fenchel Divergence* \mathcal{F} defined by

$$\mathcal{F}(x_1, x_2) := \frac{1}{2}x_1 \odot x_1 + \frac{1}{2}(x_2)_+ \odot (x_2)_+ - x_1 \odot x_2$$

We use the term *divergence* because by construction $\mathcal{F}(x_1, x_2) \geq 0$ for all $x_1, x_2 \in \mathbb{R}_+^h \times \mathbb{R}^h$.

Remark 2. We have that $\mathcal{F}(x_1, x_2) = 0$ is equivalent to $1_h^\top \mathcal{F}(x_1, x_2) = 0$. And,

$$1_h^\top \mathcal{F}(x_1, x_2) = \frac{1}{2}\|x_1\|^2 + \frac{1}{2}\|(x_2)_+\|^2 - x_1^\top x_2$$

This quantity is smooth and bi-convex in (x_1, x_2) and is equal to zero (if $x_1 \geq 0$) if and only if $x_1 = (x_2)_+$. A more natural criterion for assessing $x_1 = (x_2)_+$ would be to simply use

$$\frac{1}{2}\|x_1 - (x_2)_+\|^2 = \frac{1}{2}\|x_1\|^2 + \frac{1}{2}\|(x_2)_+\|^2 - x_1^\top (x_2)_+$$

although being seemingly close to Fenchel divergence it is easy to prove that this criterion is not differentiable and is not bi-convex in (x_1, x_2) . It is by exploiting these two essential properties of Fenchel divergence that we will be able to produce in 3.3 a bi-convex and bi-smooth formulation for implicit deep learning.

Given $X = [x_1, \dots, x_m]$ and $Z = [z_1, \dots, z_m]$ we write

$$\mathcal{F}(X, Z) = \frac{1}{m} \sum_{i=1}^m \mathcal{F}(x_i, z_i)$$

A Lagrangian relaxation approach to solving (2) using the implicit constraint formulation (3) consists in solving given a dual variable $\lambda \in \mathbb{R}_+^h$,

$$\begin{aligned} \min_{X \geq 0, \Theta} \quad & \mathcal{L}(Y, [\Theta, X]) + \lambda^\top \mathcal{F}(X, DX + EU + f1_m^\top) \\ \text{s.to:} \quad & \|D\|_2 < 1 \end{aligned} \quad (4)$$

This problem is bi-smooth in $[\Theta, X]$, but it is not convex or bi-convex. Nevertheless we can make it bi-convex with an extra conditions on Θ as shown in the next section.

3.3 Linear matrix inequality parameter constraints for bi-convexity

Let us define $\Lambda = \text{diag}(\lambda) \in \mathbb{S}_+^h$

Theorem 1. *Problem (4) is bi-convex in (Θ, X) if we impose one of the two following feasible linear matrix inequalities (LMI)*

$$\Lambda - (\Lambda D + D^\top \Lambda) \in \mathbb{S}_+^h \quad (5)$$

$$\Lambda + A^\top A - (\Lambda D + D^\top \Lambda) \in \mathbb{S}_+^h \quad (6)$$

Proof. The loss term $\mathcal{L}(Y, [\Theta, X])$ is already bi-convex in (Θ, X) , but it is not the case for the Fenchel Divergence term $\lambda^\top \mathcal{F}(X, DX + EU + f1_m^\top)$, which is not convex in X in the general case. A sufficient condition for this term to be convex in X given Θ is for the following function

$$x \rightarrow \lambda^\top \left(\frac{1}{2} x \odot x - x \odot Dx \right) = \frac{1}{2} x^\top (\Lambda - (\Lambda D + D^\top \Lambda)) x$$

to be convex. This term is convex in x if the LMI (5) is satisfied. Now the second LMI similarly arises by leveraging the fact that we can also use the term in the loss to make the objective convex in x . Indeed the objective function of (2) is convex in x if

$$x \rightarrow \frac{1}{2} x^\top A^\top A x + \frac{1}{2} x^\top (\Lambda - (\Lambda D + D^\top \Lambda)) x$$

is convex, which corresponds to LMI (6). It might not be obvious that (6) is actually an LMI, but using Schur complement we can prove it is equivalent to

$$- \begin{bmatrix} I_p & A \\ A^\top & \Lambda D + D^\top \Lambda - \Lambda \end{bmatrix} \in \mathbb{S}_+^{p+h}$$

If D satisfies (5) then it satisfies (6). We immediately have that $D = \delta I_n$ with $\delta \leq \frac{1}{2}$ satisfies (5) (and $\|D\|_2 \leq 1 - \varepsilon$). Which proves that both LMIs are feasible. \square

All in all the bi-convex problem formulation reads

$$\begin{aligned} \min_{X \geq 0, \Theta} \quad & \mathcal{L}(Y, [\Theta, X]) + \lambda^\top \mathcal{F}(X, DX + EU + f1_m^\top) \\ \text{s.to:} \quad & \|D\|_2 \leq 1 - \varepsilon, \quad \Lambda + A^\top A - (\Lambda D + D^\top \Lambda) \in \mathbb{S}_+^h, \quad \Lambda := \text{diag}(\lambda) \end{aligned} \quad (7)$$

this problem is well-posed – feasible solutions exist – and bi-smooth.

```

for  $k = 1, 2, \dots$  do
     $\Theta^k \in \operatorname{argmin}_{\Theta} \frac{1}{2m} \|AX^{k-1} + BU + c1_m^\top - Y\|_F^2 + \lambda^\top \mathcal{F}(X^{k-1}, DX^{k-1} + EU + f1_m^\top)$ 
     $\Lambda + A^\top A - (\Lambda D + D^\top \Lambda) \in \mathbb{S}_+^h$ 
     $\|D\|_2 \leq 1 - \varepsilon$ 
     $X^k \in \operatorname{argmin}_{X \geq 0} \frac{1}{2m} \|A^k X + B^k U + c^k 1_m^\top - Y\|_F^2 + \lambda^\top \mathcal{F}(X, D^k X + E^k U + f^k 1_m^\top)$ 
end

```

3.4 Block coordinate descent and first order methods

As problem (7) is bi-convex, a natural strategy is the use of block coordinate descent (BCD): alternating optimization between Θ and X . BCD corresponds to the following algorithm,

In practice such updates might be too heavy computationally as the number of datapoints m increase, or as the model size increases (i.e. h, n or p). Instead we propose to do block coordinate projected gradient updates. This method is also considered to be better at avoiding local minima. Let us denote

$$\mathcal{G}(\Theta, X) := \mathcal{L}(Y, [\Theta, X]) + \lambda^\top \mathcal{F}(X, DX + EU + f1_m^\top)$$

In the remainder of this section we derive the gradients $\nabla_{\Theta} \mathcal{G}(\Theta, X)$, $\nabla_X \mathcal{G}(\Theta, X)$ and corresponding 'optimal' step-sizes using the Lipschitz coefficients of the gradients— which is the advantage of having a bi-smooth optimization problem. Note that the objective \mathcal{G} , given X is separable in $\Theta_1 := (A, B, c)$ and $\Theta_2 := (D, E, f)$. Using scalar by matrix calculus

$$\begin{cases} \nabla_A \mathcal{G}(\Theta, X) &= \Omega(A, B, c) X^\top \in \mathbb{R}^{p \times h} \\ \nabla_B \mathcal{G}(\Theta, X) &= \Omega(A, B, c) U^\top \in \mathbb{R}^{p \times n} \\ \nabla_c \mathcal{G}(\Theta, X) &= \Omega(A, B, c) 1_m \in \mathbb{R}^p \end{cases}$$

with $\Omega(A, B, c) := \frac{1}{m} (AX + BU + c1_m^\top - Y) \in \mathbb{R}^{p \times m}$. Hence we can show that a Lipschitz constant for the gradient is given by

$$L_{\Theta_1}(X) := \frac{1}{m} \max(m, \|X\|_2^2, \|U\|_2^2, \|XU^\top\|_2)$$

the 'optimal' step-size for gradient descent is then simply given by

$$\alpha_{\Theta_1}(X) := \frac{1}{L_{\Theta_1}(X)}$$

Regarding the gradient with respect to Θ_2 , we have

$$\begin{cases} \nabla_D \mathcal{G}(\Theta, X) &= \Omega(D, E, f, \Lambda) X^\top \in \mathbb{R}^{h \times h} \\ \nabla_E \mathcal{G}(\Theta, X) &= \Omega(D, E, f, \Lambda) U^\top \in \mathbb{R}^{h \times n} \\ \nabla_f \mathcal{G}(\Theta, X) &= \Omega(D, E, f, \Lambda) 1_m \in \mathbb{R}^h \end{cases}$$

with $\Omega(D, E, f, \Lambda) := \frac{\Lambda}{m} ((DX + EU + f1_m^\top)_+ - X) \in \mathbb{R}^{h \times m}$, we can show that a Lipschitz constant for the gradient is

$$L_{\Theta_2}(X) := \frac{\lambda_{\max}}{m} \max(m, \|X\|_2^2, \|U\|_2^2, \|X\|_2 \|U\|_2)$$

where $\lambda_{\max} = \max_{j \in \{1, \dots, h\}} \lambda_j$. We can then similarly define an 'optimal' step-size α_{Θ_2} . We have that

$$\nabla_X \mathcal{G}(\Theta, X) = \frac{1}{m} \left\{ A^\top (AX + BU + c1_m^\top) + (\Lambda - \Lambda D - D^\top \Lambda) X + D^\top \Lambda (DX + EU + f1_m^\top)_+ - \Lambda (EU + f1_m^\top) \right\}$$

A Lipschitz constant for this gradient is

$$L_X(\Theta) = \frac{1}{m} (\|A^\top A + \Lambda - \Lambda D + D^\top \Lambda\|_2 + \lambda_{\max} \|D\|_2^2)$$

we can then take the step-size $\alpha_X(\Theta) = 1/L_X(\Theta)$. We propose the following block coordinate projected gradient scheme (BC-gradient) to find a candidate solution to (7). We denote compactly the convex set

$$\mathcal{S}_\Theta := \{\Theta \mid \Lambda + A^\top A - (\Lambda D + D^\top \Lambda) \in \mathbb{S}_+^h, \quad \|D\|_2 \leq 1 - \varepsilon\}$$

and $\mathcal{P}_{\mathcal{S}_\Theta}$ the corresponding convex projection

```

for  $k = 1, \dots$  do
     $\Theta^k = \mathcal{P}_{\mathcal{S}_\Theta} \left( \Theta^{k-1} - \alpha_\Theta(X^{k-1}) \nabla_\Theta \mathcal{G}(\Theta^{k-1}, X^{k-1}) \right)$ 
     $X^k = \left( X^{k-1} - \alpha_X(\Theta^k) \nabla_X \mathcal{G}(\Theta^k, X^{k-1}) \right)_+$ 
end

```

3.5 Dual methods

We propose the following schemes to find an appropriate dual variable λ . Let $\varepsilon > 0$ be a precision parameter for the implicit constraint, I.e. such that we would have

$$\mathcal{F}(X, DX + EU + f1_m^\top) \leq \varepsilon$$

We start with $\lambda = 0$ and we solve the two following separate problems

$$\min_{X \geq 0, A, B, c} \frac{1}{m} \|AX + BU + c1_m^\top - Y\|_F^2$$

and then

$$\min_{D, E, f} 1_h^\top \mathcal{F}(X, DX + EU + f1_m^\top)$$

If $\mathcal{F}^* := \mathcal{F}(X, DX + EU + f1_m^\top) < \varepsilon I_h$ then we stop there. Otherwise, we do one of the two following 'dual updates'

3.5.1 Dual ascent conditional on Fenchel Divergence

$$\lambda \leftarrow \lambda + \alpha \mathcal{F}^* \odot 1\{\mathcal{F}^* \geq \varepsilon I_h\} \quad (8)$$

where $\alpha > 0$ is a step-size. Note that here we only update the components of λ for which the corresponding Fenchel divergence is more than ε . We then proceed to solve (7) using previously discussed methods and iterate. Alternatively, if the BC-gradient method is used, we can do a dual update after each BC-gradient update.

3.5.2 Dual variable update conditional on loss

We start with $\lambda = \varepsilon I_h$. Given (Θ, X) , we define the unique \overline{X} such that the implicit constraint is enforced given Θ

$$\overline{X} = (D\overline{X} + EU + f1_m^\top)_+$$

We then define $\Delta X := X - \overline{X}$. We can compute in close form the error on the loss due to the implicit constraint violation

$$\begin{aligned} \Delta \mathcal{L} &:= \mathcal{L}(Y, [\Theta, \overline{X}]) - \mathcal{L}(Y, [\Theta, X]) \\ &= \frac{1}{2m} \left(\|A\Delta X\|_F^2 + \text{Tr}(\Omega, A\Delta X) \right) \end{aligned}$$

with $\Omega := BU + c1_m^\top$. We can write this error as a sum of contributions with respect to each hidden variable components $j \in \{1, \dots, h\}$

$$\Delta \mathcal{L} = \sum_{j=1}^h \left\{ \Delta \mathcal{L}_j := \frac{1}{m} A_j^\top \left(\frac{1}{2} A\Delta X + \Omega \right) \Delta X_j^\top \right\}$$

where $A_j \in \mathbb{R}^h$ is the j -th column of A and $\Delta X_j \in \mathbb{R}^{1 \times m}$ is the j -th row of ΔX . The objective of this dual update is to achieve an error on the loss that is smaller than a fraction $\eta \in (0, 1)$ of the loss

$$\frac{\Delta \mathcal{L}}{\mathcal{L}(Y, [\Theta, \bar{X}])} \leq \eta$$

In order to update each component of the dual variable, we propose the following update. Given $j \in \{1, \dots, h\}$ if

$$\frac{(\Delta \mathcal{L}_j)_+}{\mathcal{L}(Y, [\Theta, \bar{X}])} \geq \frac{\eta}{h}$$

then we do the update

$$\lambda_j \rightarrow \beta \lambda_j$$

with $\beta > 1$ a hyperparameter.

References

- [1] L. E. Ghaoui, F. Gu, B. Travacca, and A. Askari. Implicit deep learning, 2019.